# PROJECT NASTY- NEWS AGGREGATION AND SUMMARIZATION

## A

## MAJOR PROJECT REPORT

**Submitted by**
**SUPRATIK CHATTERJEE (RA1511003040196)**

*Under the guidance of*

## Mr. Muthurasu N., B.E., M.Tech.,

*(Assistant professor, Department of Computer Science Engineering)*

*in partial fulfilment of the degree*
*Of*

## BACHELOR OF TECHNOLOGY

*in*

## COMPUTER SCIENCE ENGINEERING

*of*

## FACULTY OF ENGINEERING AND TECHNOLOGY



SRM INSTITUTE OF TECHNOLOGY AND SCIENCE
Vadapalani, Chennai-600026

**APRIL 2019**

# PROJECT NASTY- NEWS AGGREGATION AND SUMMARIZATION

## A

## MAJOR PROJECT REPORT

**Submitted by**
**SUPRATIK CHATTERJEE (RA1511003040196)**

*Under the guidance of*

## Mr. Muthurasu N., B.E., M.Tech.,

*(Assistant professor, Department of Computer Science Engineering)*

*in partial fulfilment of the degree*
*Of*

## BACHELOR OF TECHNOLOGY

*in*

## COMPUTER SCIENCE ENGINEERING

*of*

## FACULTY OF ENGINEERING AND TECHNOLOGY



SRM INSTITUTE OF TECHNOLOGY AND SCIENCE
Vadapalani, Chennai-600026

**APRIL 2019**

# BONAFIDE CERTIFICATE

It is certified that this project report titled "Project NASTY" is the bonafide work of **Supratik Chatterjee (RA1511003040196)** , who carried out the project work under my supervision.

**Signature of the Guide**                    **Signature of the HOD**

Mr. N. Muthurasu                              Dr. S. Prasanna Devi
B.E., M.Tech,                                 B.E., M.Tech, Ph.D.,
                                              PGDHRM, PDF(IISc),
Assistant Professor                           Head of the Department
Department of CSE                             Department of CSE
SRM IST                                       SRM IST
Vadapalani,Chennai-                           Vadapalani,Chennai-
600026                                        600026

# Acknowledgement

I place my gratitude with our respected chancellor of SRM Institute of Science and Technology, to have provided the required infrastructure, to hone my skills. I also thank my teachers for their guidance and teachings, without which I would not be able to learn what I know as efficiently as I do.

We extend this opportunity to extend our heartfelt gratitude to our respected Dean **Dr. K. Duraivelu** for his impeccable guidance.

We deeply express our sincere thanks to our Head of Department **Dr. S. Prasanna Devi**, for encouraging and allowing us to present the project on the topic, **News Analysis and Summarization Tool**, at our department premises for the partial fulfillment of the requirements leading to the award of B-Tech degree.

It is our privilege to express our sincerest regards to our project coordinator, **Mr. Muthurasu N.** for his valuable inputs, able guidance, encouragement, whole-hearted cooperation and constructive criticism throughout the duration of our project.

We take this opportunity to thank all our faculty members and management who have directly or indirectly helped our project. Last but not the least we express our thanks to our friends for their cooperation and support.

# Abstract

In the modern day field of media, news is majorly affected by the social platforms such as Facebook, Twitter, etc. It becomes necessary to understand all the facts of an event, and to quantify, or relate them to understand it. Most humans do not have the dedication or time to study all the events at all the time, as time is always limited. However, an aggregation and analysis system on a computer could do it for them. A computer can form links and store endless amounts of data which would only be affected by the amount of storage space available and the storage algorithm being made use of. Apart from that advantage we have the advantage of having programs such as Hadoop and MongoDB to take care of a lot of work for us. This project has been worked upon with an idea in mind, to create a common platform for people to fetch aggregated data from multiple news agencies through an abstraction.

# Contents

# List of Figures

# CHAPTER 1
## Introduction

When dealing with news papers, it becomes a tedious chore to understand which one to follow and why. A lot of times, we see that in order for people to report an event, they sometimes have to withhold some little pieces only to reveal those at a later date. Apart from these, there is also the part where they need to be stored categorically to create a record of the modern history which keeps growing more and more opaque due to the complications everyday. So, how do we keep track of all facts while making sure we understand all related facets to an event?

The solution is by relying on autonomous systems to work towards continuous content fetching, analyzing, aggregating and relating facts. But where can a person find the facts? It so happens that most major news channels offer a format known as RSS, to fetch the data. However, it is nigh impossible for us to keep track with the humongous number of news networks that exist within our country alone (the number being 1,05,443, of which, 13,661 are English Newspapers), let alone that of news channels all over the world. So in order to reduce the generation of noisy data, we reduce the scope to work with the most prominents ones.

The prominent news channels provide a certain utility known as an Atom feed, or RSS feed. These feeds have been shown in this section. These utilities are however under-utilized in most cases. Our work is to exploit this resource, made available from news networks, and manage it to maximize their usage.
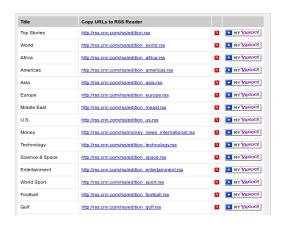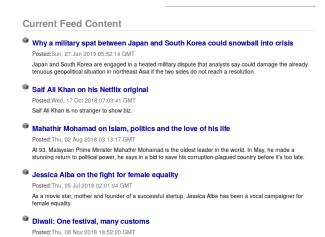
Figure 1.1: CNN RSS Listings



Figure 1.2: CNN RSS feed



Figure 1.3: BBC RSS Listings



Figure 1.4: BBC RSS feed

# CHAPTER 2
# Project Work

Facts collection, aggregation and relation is everything this project is about. RSS feeds usually contain the shortened form of all the events, covering only the facts of the event.

## 2.1 The Concept

News articles contain well structured sentences that contain information in a human-readable language using the least amount of complicated words as can be managed.

To obtain the content, we make use of a Focused Incremental Web Crawler(A web crawler that extracts only specific content at a relatively shallow depth). This content is brought and stored for consideration in a pre-determined storage unit. This then moves on for analysis.

Facts aggregation is a lot about analysis of the text in the content, to find for three primary things, which are attributes, relations and involvement in events. These three are the specific information required. This involves the usage of NLP methodologies such as Named Entity Relations, Parts of Speech Tagging, Anaphora Resolution, and several others, which are to work together to reduce the amount of useless data to be stored about entities.

## 2.2 Problem Statement

Aggregation and analysis of news over network, from several channels to remove repetitions and increase efficient centralization of factual data of any event, location or an object.

## 2.3 Proposed Solution

The solution is to make use of an observer program to do it all, at regular intervals by utilizing a programming language, which in this particular case, for prototyping, would be through python.

## 2.4 Obstacles

There are three obstacles to this:

### 2.4.1 Relation Extraction

### 2.4.2 Anaphora Resolution

### 2.4.3 Metadata Repository

## 2.5 Advantages

## 2.6 Future Scope

## 2.7 Available tools

There are absolutely no tools available of this kind. There are several that forward the RSS feeds onto devices. But none that analyze to obtain all facts and aggregate it at a single point. A few of the more remarkable ones are specified here for your consideration.

### 2.7.1 News API

## 2.8 Tools Studied and Used

### 2.8.1 Python 3

### 2.8.2 nltk

### 2.8.3 Stanford CoreNLP

### 2.8.4 RSS

### 2.8.5 JSON

### 2.8.6 MongoDB

## 2.9 Theories

# CHAPTER 3
## Conclusion